



Munich Personal RePEc Archive

# **ivporbit:An R package to estimate the probit model with continuous endogenous regressors**

Taha Zaghdoudi

University of Jendouba, Faculty of Law Economics and Management  
of Jendouba

25 September 2014

Online at <https://mpra.ub.uni-muenchen.de/72383/>  
MPRA Paper No. 72383, posted 6 July 2016 06:59 UTC

# ivprobit: An R package to estimate the probit model with continuous endogenous regressors

Taha Zaghdoudi\*

25-10-2014

## Abstract

One of the most important problem of misspecification in the probit model is the correlation between regressors and error term. To deal with this problem, some commercial software gives a solution such as Stata. For the famous R language the `ivprobit` gives the users the way to estimate the instrumental probit model.

## 1 Introduction

Many econometrics study are focused in various kinds of misspecification in the limited-dependent variable models such as correlation between regressors and error term which produce inconsistent results. To avoid this problem we apply the instrumental variable method in which we use the correlated variables as instruments. A two-stage method are used by some authors [Blundell and Powell \[2004\]](#) to fit the probit model but it produce non efficient result.

Otherwise, [Newey \[1987\]](#) expose an efficient way to estimate limited-dependent variable model by using the Amemiya's Generalized Least Squares estimators [Amemiya \[1978\]](#). The idea is to include in the two-stage model a continuous endogenous regressor. This method is used in Stata when the MLE fail to estimate the model. For the `ivprobit` we apply the same method as used in Stata.

In the following sections we expose the Newey method to estimate the instrumental probit model then we give a simple example how to use the `ivprobit` Package.

## 2 Estimation method

Generally the model is:

$$\begin{aligned}y_{1i}^* &= y_{2i}\beta + x_{1i}\gamma + \mu_i \\ y_{2i} &= x_{1i}\Pi_1 + x_{2i}\Pi_2 + v_i.\end{aligned}$$

---

\*University of Jendouba, Faculty of Law, Economics and Management of Jendouba

Where  $i = 1, \dots, N$ ,  $y_{2i}$  is a  $1 \times p$  vector of endogenous variables,  $x_{1i}$  is a  $1 \times 1$  vector of exogenous variable,  $x_{2i}$  is a  $1 \times k_2$  vector of additional instruments, and the equation for  $y_{2i}$  is written in reduced form. By assumption,  $(v_i, \mu_i) \sim N(0, \Sigma)$ , where  $\sigma_{11}$  is normalized to one to identify the model.  $\beta$  and  $\gamma$  are vectors of structural parameters, and  $\Pi_1$  and  $\Pi_2$  are matrices of reduced-form parameters. This is a recursive model:  $y_{2i}$  appears in the equation for  $y_{1i}^*$ , but  $y_{1i}^*$  does not appear in the equation for  $y_{2i}$ . We do not observe  $y_{1i}^*$ ; instead, we observe:

$$y_{1i} = \begin{cases} 0 & y_{1i}^* \geq 0 \\ 1 & y_{1i}^* \leq 1 \end{cases}$$

The order condition for identification of the structural parameters requires that  $k_2 \geq 1$ . Presumably, is not block diagonal between  $\mu_i$  and  $v_i$ ; otherwise,  $y_{2i}$  would not be endogenous.

To obtain the two-step estimates, Newey [1987] use minimum chi-squared estimator and the model is:

$$y_{1i}^* = z_i \delta + \mu_i \quad y_{2i} = x_i \Pi + v_i \quad (1a)$$

where  $z_i = (y_{2i}, x_{1i})$ ,  $x_i = (x_{1i}, x_{2i})$ ,  $\delta = (\beta', \delta')'$ , and  $\Pi = (\Pi_1', \Pi_2')'$ . (1b)

The reduced-form equation for  $y_{1i}^*$  is:

$$y_{1i}^* = (x_i \Pi + v_i) \beta + x_{1i} \gamma + \mu_i = x_i \alpha + v_i \beta + \mu_i = x_i \alpha + \nu_i$$

where

$\nu_i = v_i \beta + \mu_i$ . Because  $\mu_i$  and  $v_i$  are jointly normal,  $\nu_i$  is also normal. Note that:

$$\alpha = \{\Pi_1\} \beta + \{\gamma\} = D(\Pi) \delta$$

where  $D = (\Pi) = (\Pi, I_1)$  and  $I_1$  is defined such that  $x_i I_1$ . Letting  $\hat{z}_i = (x_i \hat{\Pi}, x_{1i})$ ,  $\hat{z}_i = x_i D(\hat{\Pi}) \delta$ , where  $D(\hat{\Pi}) = D(\hat{\Pi}, I_1)$ . Thus one estimator of  $\alpha$  is  $D(\hat{\Pi}) \delta$ ; denote this estimator by  $\hat{D} \delta$ .

$\alpha$  could also be estimated directly as the solution to:

$$\max_{\alpha, \lambda} = \sum_{i=1}^N l(y_{1i}, x_i \alpha + \hat{v}_i \lambda)$$

where  $l(\cdot)$  is the log likelihood for probit. Denote this estimator by  $\hat{\alpha}$ . The inclusion of the term  $\hat{v}_i \lambda$  follows because the multivariate normality of  $(\mu_i, v_i)$  implies that, conditional on  $y_{2i}$ , the expected value of  $\mu_i$  is nonzero. Because  $v_i$  is unobservable, the least-squares residuals from fitting (1b) are used.

Amemiya [1978] shows that the estimator of  $\delta$ :

$$\max_{\delta} (\tilde{\alpha} - \hat{D} \delta)' \hat{\Omega}^{-1} (\tilde{\alpha} - \hat{D} \delta) \quad (2)$$

where  $\hat{\Omega}$  is a consistent estimator of the covariance of  $\sqrt{N}(\tilde{\alpha} - \hat{D} \delta)$ , is asymptotically efficient relative to all other estimators that minimize the distance between  $\tilde{\alpha}$  and  $D(\hat{\Pi}) \delta$ . Thus an efficient estimator of  $\delta$  is :

$$\hat{\delta} = (\hat{D}' \hat{\Omega}^{-1} \hat{D})^{-1} (\hat{D}' \hat{\Omega}^{-1} \tilde{\alpha}) \quad (3)$$

$$\text{and} \\ \text{Var}(\hat{\delta}) = (\hat{D}' \hat{\Omega}^{-1} \hat{D})^{-1} \quad (4)$$

To implement this estimator, we need  $\hat{\Omega}^{-1}$ . Consider the two-step maximum likelihood estimator that results from first fitting (1b) by OLS and computing the residuals  $\hat{v}_i = y_{2i} - x_i' \hat{\Pi}$ . The estimator is then obtained by solving:

$$\max_{\delta, \lambda} = \sum_{i=1}^N l(y_{1i}, z_i \delta + \hat{v}_i \lambda)$$

This is the two-step instrumental variables (2SIV) estimator proposed by [Rivers and Vuong \[1988\]](#), and its role will become apparent shortly.

From Proposition 5 of [\[Newey, 1987\]](#),  $\sqrt{N}(\tilde{\alpha} - \hat{D}\delta) \xrightarrow{d} N(0, \Omega)$ , where

$$\Omega = J_{\alpha\alpha}^{-1} + (\lambda - \beta)' \Sigma_{22} (\lambda - \beta) Q^{-1}$$

and  $\Sigma_{22} = E\{v_i' v_i\} \cdot J_{\alpha\alpha}^{-1}$  is simply the covariance matrix of  $\tilde{\alpha}$ , ignoring that  $\hat{\Pi}$  is an estimated parameter matrix. Moreover, Newey shows that the covariance matrix from an OLS regression of  $y_{i2}(\hat{\lambda} - \hat{\beta})$  on  $x_i$  is a consistent estimator of the second term.  $\hat{\lambda}$  can be obtained from solving (2), and the 2SIV estimator yields a consistent estimate,  $\hat{\beta}$ .

The basic algorithm proceeds as follows:

1. Each of the endogenous right-hand-side variables is regressed on all the exogenous variables, and the fitted values and residuals are calculated. The matrix  $\hat{D} = D(\hat{\Pi})$  is assembled from the estimated coefficients.
2. probit is used to solve (2) and obtain  $\tilde{\alpha}$  and  $\tilde{\lambda}$ . The portion of the covariance matrix corresponding to  $\alpha$ ,  $J_{\alpha\alpha}^{-1}$ , is also saved.
3. The 2SIV estimator is evaluated, and the parameters  $\hat{\beta}$  corresponding to  $y_{2i}$  are collected.
4.  $y_{i2}(\hat{\lambda} - \hat{\beta})$  is regressed on  $x_i$ . The covariance matrix of the parameters from this regression is added to  $J_{\alpha\alpha}^{-1}$ , yielding  $\hat{\Omega}$ .
5. Evaluating (3) and (4) yields the estimates  $\hat{\delta}$  and  $\text{Var}(\hat{\delta})$ .
6. A Wald test of the null hypothesis  $H_0 : \lambda=0$ , using the 2SIV estimates, serves as our test of exogeneity.

### 3 Example

To download the package you can just type `install.packages("ivprobit")` in R or Rstudio.

```
# load the package
library(ivprobit)
```

The data we use in this example represents the Foreign-Exchange Derivatives Use By Large U.S. Bank Holding Companies from 1996 to 2000.

```
# load the database
dat<-system.file("data", "eco.RData", package="ivprobit")
load(dat)
pro<-library("ivprobit", lib.loc="C:/Program Files/R/R-3.0.3/library")
#load data
dat<-system.file("data", "eco.RData", package="ivprobit")
load(dat)
```

The function is `ivprobit(y x,y2 z, data)`:

- `y`: the dichotomous l.h.s vector
- `x`: the r.h.s. exogenous variables matrix.
- `y2`: the r.h.s. endogenous variables vector or matrix
- `z`: the complete set of instruments (a matrix)
- `data`: the dataframe

In this example we have `d2` the dichotomous l.h.s vector, the r.h.s exogenous variables are (`ltass,roe` and `div`), the r.h.s. endogenous variables matrix is (`eqrat,bonus`) and the instrumental variables are (`ltass,roe,div,gap,cfa`).

```
# fit the instrumental probit model
pro<-ivprobit(d2~ltass+roe+div,cbind(eqrat,bonus)~ltass+roe+div+gap+cfa,mydata)
# the results summary
summary(pro)
```

##		Coef	S.E.	t-stat	p-val
##	(Intercept)	-1.6862e+01	9.7317e+00	-1.7327	0.08355 .
##	X1ltass	9.4760e-01	6.2070e-01	1.5267	0.12724
##	X1roe	6.6492e-02	1.1312e-01	0.5878	0.55684
##	X1div	2.0378e-06	4.2745e-06	0.4767	0.63368
##	yhateqrat	9.4371e+00	1.2937e+01	0.7295	0.46593
##	yhatbonus	-1.4173e-06	2.8163e-06	-0.5032	0.61493
##	---				
##	Signif. codes:	0 '***'	0.001 '**'	0.01 '*'	0.05 '.' 0.1 ' ' 1

## 4 Conclusion

In this paper we introduced the `ivprobit` to estimate the parameters of an dichotomous choice model that contains endogenous regressors. The routine is simple and yields the same results as the two-step option in the commercially available Stata software.

## References

- Takeshi Amemiya. The estimation of a simultaneous equation generalized probit model. *Econometrica: Journal of the Econometric Society*, 46:1193–1205, 1978.
- Richard W Blundell and James L Powell. Endogeneity in semiparametric binary response models. *The Review of Economic Studies*, 71(3):655–679, 2004.
- Whitney K Newey. Efficient estimation of limited dependent variable models with endogenous explanatory variables. *Journal of Econometrics*, 36(3):231–250, 1987.
- Douglas Rivers and Quang H Vuong. Limited information estimators and exogeneity tests for simultaneous probit models. *Journal of Econometrics*, 39(3):347–366, 1988.